

```

1  --
2
3
4
5
6
7
8
9 local l = require"lhb"
10 local the = l.settings[[
11   rl.lua : strings
12   (c)2022 Tim Menzies <tim@ieee.org> BSD-2clause.
13
14 OPTIONS:
15 -b --bins discretization control = 8
16 -k --keep keep only these nums = 256]]
17
18 local About = {} -- factor for making columns
19 local Col = {} -- summarize one column
20 local Data = {} -- store rows, and their column summaries
21 local Row = {} -- store one row
22
23 -- CODE CONVENTIONS
24 -- Leading__upper_case : class
25 -- l. : instance va
26 -- l. s : reference to a library function
27 -- prefix _ : some internal function,variable.
28
29 -- type hints: where practical, on function arguments,
30 -- t = table
31 -- _ prefix s=string
32 -- _ prefix n=num
33 -- _ prefix is=boolean
34 -- _ class names in lower case denote vars of that class
35 -- _ suffix s denotes table of things
36
37 -----
38
39 -- About
40
41
42 function About.new(sNames)
43   return About._cols((names=sNames, all={}, x={}, y={}, klass=nil),sNames) end
44
45 -- How to recognize different column types
46 local _is={
47   nom = "a-z", -- ratio cols start with uppercase
48   goal = "[!a-z]", -- !klass, {t, _}=maximize,minimize
49   klass = "S", -- klass if "!"
50   skip = "$", -- skip if ":"
51   less = "$", -- minimize if "-"
52 }
53
54 -- Turn a list of column names into Col objects. If the new col is independent
55 -- or dependent or a goal attribute then remember that in i.x or i.y or i.klass.
56 function About._cols(i,sNames)
57   for at,name in pairs(sNames) do
58     local col = l.push(i.all, Col.new(name,at))
59     if not name:find(_is.skip) then
60       l.push(name:find(_is.goal) and i.y or i.x, col)
61       if name:find(_is.klass) then i.klass=col end end
62     return i end
63
64 -- Update, only the non-skipped cols (i.e. those found in i.x and j.x.
65 function About.add(i,t)
66   local row = t.cells and t or Row.new(i.about, t)
67   for _,cols in pairs(i.x,i.y) do
68     for _,col in pairs(cols) do
69       Col.add(col, row.cells[col.at]) end end
70   return row end
71
72 -----
73
74 -- Col
75
76 -- Summarize one column.
77 function Col.new(txt,at)
78   txt = txt or ""
79   return {n = 0, -- how many items seen?
80     at = at or 0, -- position of column
81     txt = txt, -- column header
82     isNom = txt:find(_is.nom),
83     w = txt:find(_is.less) and -1 or 1,
84     ok = true, -- false if some update needed
85     _has = {} end -- place to keep (some) column values.
86
87 -- Update
88 function Col.add(i,x)
89   if x == "" then
90     i.n = i.n + 1
91     if i.isNom
92     then i._has[x] = 1 + (i._has[x] or 0)
93     else local pos
94       if #i._has < the.keep then pos= 1 + (#i._has)
95       elseif l.rand() < the.keep/i.n then pos=l.rand(#i._has) end
96       if pos then
97         i.ok=false -- kept items are no longer sorted
98         i._has[pos]=x end end end
99
100 -- Distance
101 function Col.dist(i,x,y)
102   if x=="?" and y=="?" then return 1 end
103   if i.isNom
104   then return x==y and 0 or 1
105   else if x=="?" and y=="?" then return 1 end
106   if x=="?" then y = Col.norm(i,y); x=y<.5 and 1 or 0
107   elseif y=="?" then x = Col.norm(i,x); y=x<.5 and 1 or 0
108   else x,y = Col.norm(i,x), Col.norm(i,y) end
109   return math.abs(x-y) end end
110
111 -- Diversity
112 function Col.div(i)
113   if i.isNom
114   then local e=0
115     for v in pairs(i._has) do
116       if v>0 then e=e-v/i.n*math.log(v/i.n,2) end end
117     return e
118   else local t=Col.has(i)
119     return (l.per(t,.9) - 1.per(t,.1))/2.56 end end
120
121 -- Sorted contents
122 function Col.has(i)
123   if i.isNom then return i._has end
124   if not i.ok then table.sort(i._has) end
125   i.ok=true

```

```

126   return i._has end
127
128 -- Central tendency
129 function Col.mid(i)
130   if i.isNom
131   then local mode,most=nil,-1
132     for k,v in pairs(i._has) do if v>most then mode,most=k,v end end
133     return mode
134   else return l.per(Col.has(i),.5) end end
135
136 -- Return num, scaled to 0..1 for lo..hi
137 function Col.norm(i,num)
138   local a = Col.has(i) -- "a" contains all our numbers, sorted.
139   return a[#a] - a[1] < 1E-9 and 0 or (num-a[1])/(a[#a]-a[1]) end
140
141 -- Map x to a small range of values.
142 function Col.discretize(i,x, a,b,lo,hi)
143   if i.isNom then return x else
144     a = has(i)
145     lo,hi = a[1], a[#a]
146     b = (hi - lo)/the.bins
147     return hi=lo and 1 or math.floor(x/b+.5)*b end end
148
149 -----
150
151 -- Row
152
153 -- Hold one record
154 function Row.new(about,t)
155   return {about=about, cells=t, cooked=l.map(t,l.same)} end
156
157 -- Everything in rows, sorted by distance to i.
158 function Row.around(i,rows)
159   local fun = function(j) return (row=j, d=Row.dist(i,j)) end
160   return l.sort(l.map(rows, fun), lt"d") end
161
162 -- Recommend sorting i before j (since i is better).
163 function Row.better(i,j)
164   i.evaled,j.evaled= true,true
165   local s1,s2,d,n,x,y=0,0,0,0
166   local ys,e = i._about.y,math.exp(1)
167   for _,col in pairs(ys) do
168     x,y = i.cells[col.at], j.cells[col.at]
169     x,y = Col.norm(col,x), Col.norm(col,y)
170     s1 = s1 - e^(Col.w * (x-y)/#ys)
171     s2 = s2 - e^(Col.w * (y-x)/#ys) end
172   return s1/#ys < s2/#ys end
173
174 -- Distance
175 function Row.dist(i,j)
176   local d,n,x,y,dist=0,0
177   local cols = cols or i._about.x
178   for _,col in pairs(cols) do
179     x,y = i.cells[col.at], j.cells[col.at]
180     d = d + Col.dist(col,x,y)*the.p
181     n = n + 1 end
182   return (d/n)^(1/the.p) end
183
184 -----
185
186 -- Data
187
188 -- Holds n records
189 function Data.new(t) return {rows={}, about=About.new(t) } end
190
191 -- Update
192 function Data.add(i,t) l.push(i.rows, About.add(i.about,t)) end
193
194 -- Load from file
195 function Data.load(sFilename, data)
196   l.csv(sFilename, function(row)
197     if data then Data.add(data,row) else data=Data.new(row) end end)
198   return data end
199
200 -- Central tendency
201 function Data.mid(i) return l.map(i.about.y, Col.mid) end
202
203 -- Discretize all row values (writing those vals to "cooked").
204 function Data.discretize(i)
205   for _,col in pairs(i.about.x) do
206     for _,row in pairs(i.rows) do
207       local x = row.cells[col.at]
208       if x == "" then
209         row.cooked[col.at] = discretize(col,x) end end end end
210
211 -----
212
213 -- M:it
214
215 d=Data.load("./Data/auto93.csv")
216
217 1.chat(Data.mid(d))
218 1.chat(d.about.x[1])
219 --map(d.about.x, chat)
220 -- chat(d.About.x)
221 -- bins(d)
222 -- for row in pairs(d.rows) do 1.chat(Row.cooked) end
223 -- for i=1,20 do
224 --   r1=l.any(d.rows)
225 --   r2=l.any(d.rows)
226 -- end

```